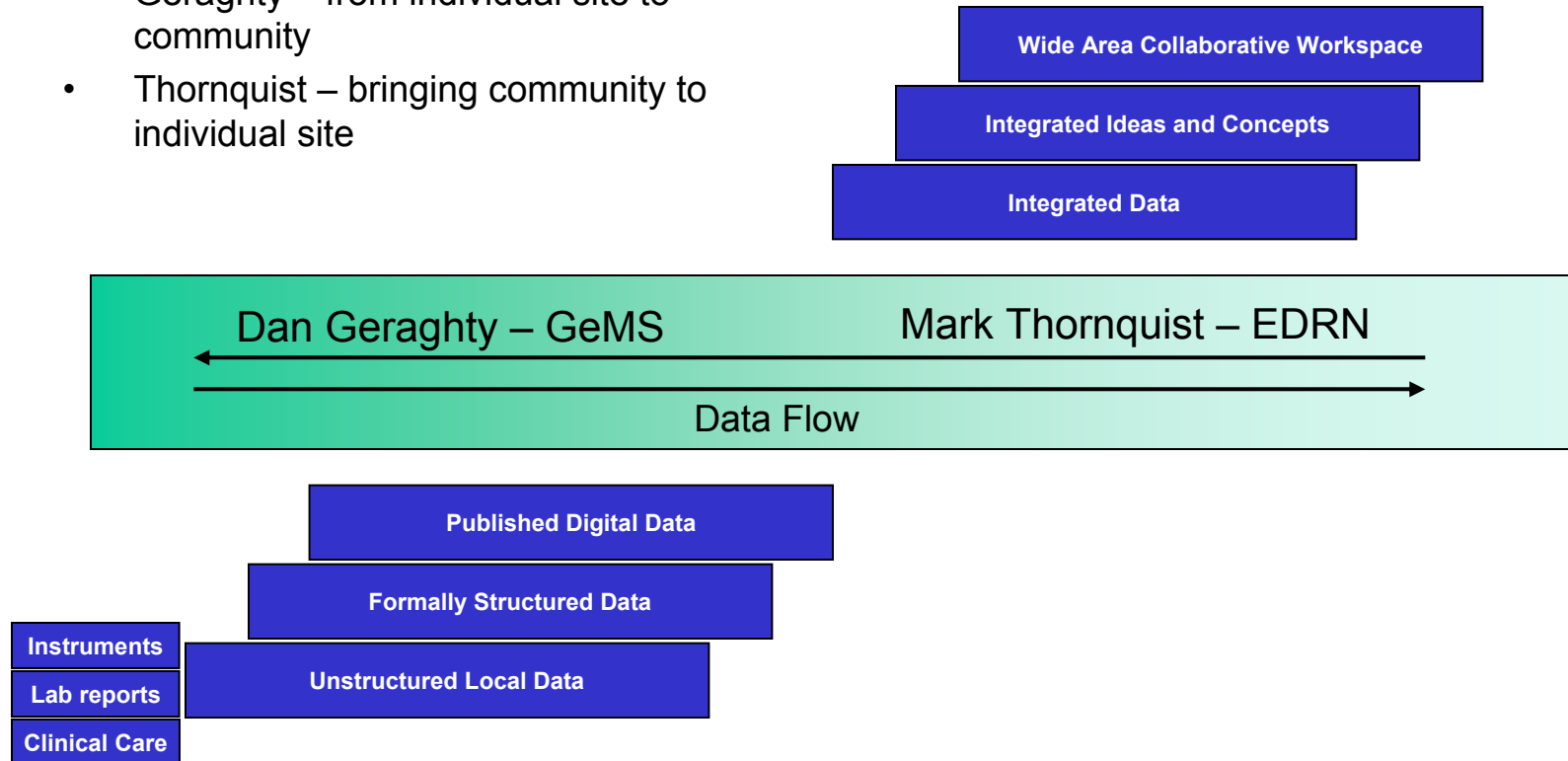# caBIG Architecture Kickoff Meeting Presentation
# Fred Hutchinson Cancer Research Center

**Mark Thornquist, Derek Walker, Heather Kincaid, Rahul Joshi, Dan Geraghty, Robert Robbins.**

# Data Sharing Continuum

- Geraghty – from individual site to community
- Thornquist – bringing community to individual site

**Wide Area Collaborative Workspace**

**Integrated Ideas and Concepts**

**Integrated Data**

Dan Geraghty – GeMS          Mark Thornquist – EDRN

Data Flow

**Published Digital Data**

**Formally Structured Data**

**Instruments**

**Lab reports**

**Unstructured Local Data**

**Clinical Care**

## Development Principles

- <u>Roadmap Driven</u>: all pieces align with a reference architecture / roadmap
- <u>Flexibility in inputs and outputs</u>: allows variety of data types and meta data classifications to  co-exist within the same system
- <u>Scalable Design</u>: retain system performance under increasing system load
- <u>Wide Ranging</u>: retain consistency with other information technology initiatives
- <u>Technology Agnostic</u>: allow for variety of technologies to exchange data
- <u>Open source</u>: allow interested parties to adopt, modify and improve the current state

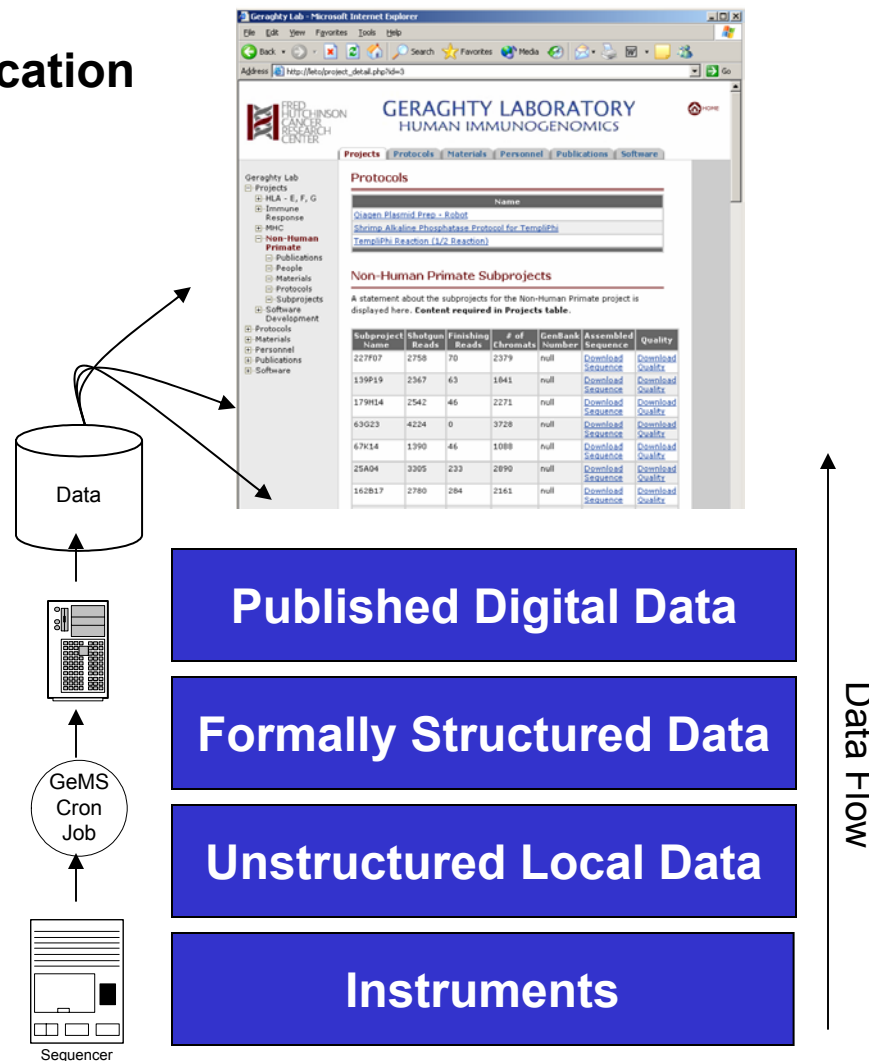# Different Approaches for Different Circumstances

## Geraghty – GeMS

- Integration through usage
- Provide useful, needed tools – resulting in *de facto* common data

## Thornquist – EDRN

- Integration through middleware
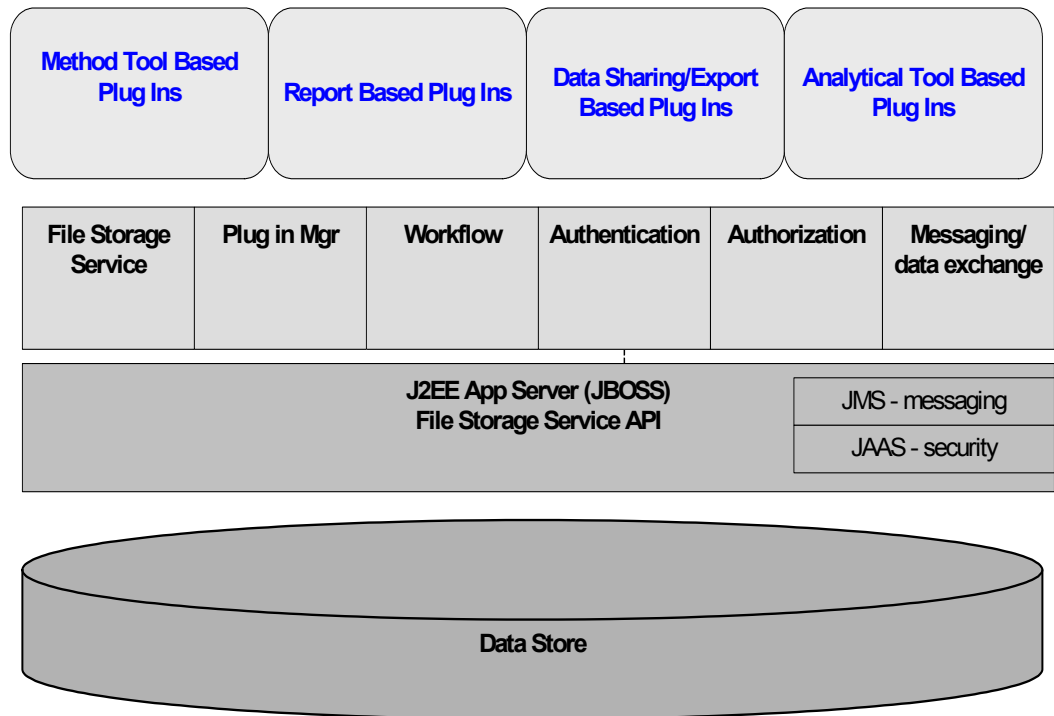- Map existing databases to common data elements

# From Data Generation to Data Publication



- Nightly Data pick up by system

- Unstructured and unrelated data sent to GeMS server for processing

- Data related to associated parameters

- Subset of data made available to the Geraghty website

**Published Digital Data**

**Formally Structured Data**

**Unstructured Local Data**

**Instruments**

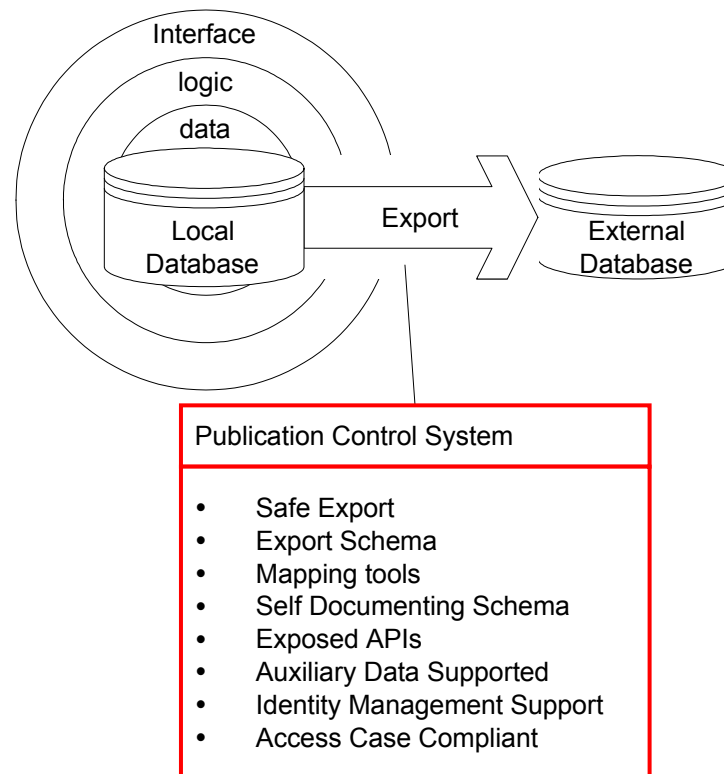Data Flow

Data

GeMS Cron Job

Sequencer

# GeMS Architecture

- The data store is accessed through a file storage service API that acts as a DAO (Data Access Object) Layer.
- Core services is made available above J2EE application server. These services are used by the plugins to carry out their functions.
  - *File Storage Service– manages file system*
  - *Authentication – identify validation*
  - *Authorization – users level of access*
  - *Messaging – local workflow processes andcollaboration with remot sites*
  - *Plugin Manager – manages the resigration of plugin components*
  - *Workflow – manages the workflow agents, their states, and the associated triggers*
- Plugins represent the functional components that use the core services.

| Method Tool Based Plug Ins | Report Based Plug Ins | Data Sharing/Export Based Plug Ins | Analytical Tool Based Plug Ins |
|---|---|---|---|

| File Storage Service | Plug in Mgr | Workflow | Authentication | Authorization | Messaging/ data exchange |
|---|---|---|---|---|---|

**J2EE App Server (JBOSS)**
**File Storage Service API**

JMS - messaging

JAAS - security

**Data Store**

# Generalizing the Data over a grid

- Next phase to build data sharing mechanism based on development of generic publication control system (export server)

- Test publication control and data sharing across disciplines with the Thornquist's EDRN/ERNE development efforts
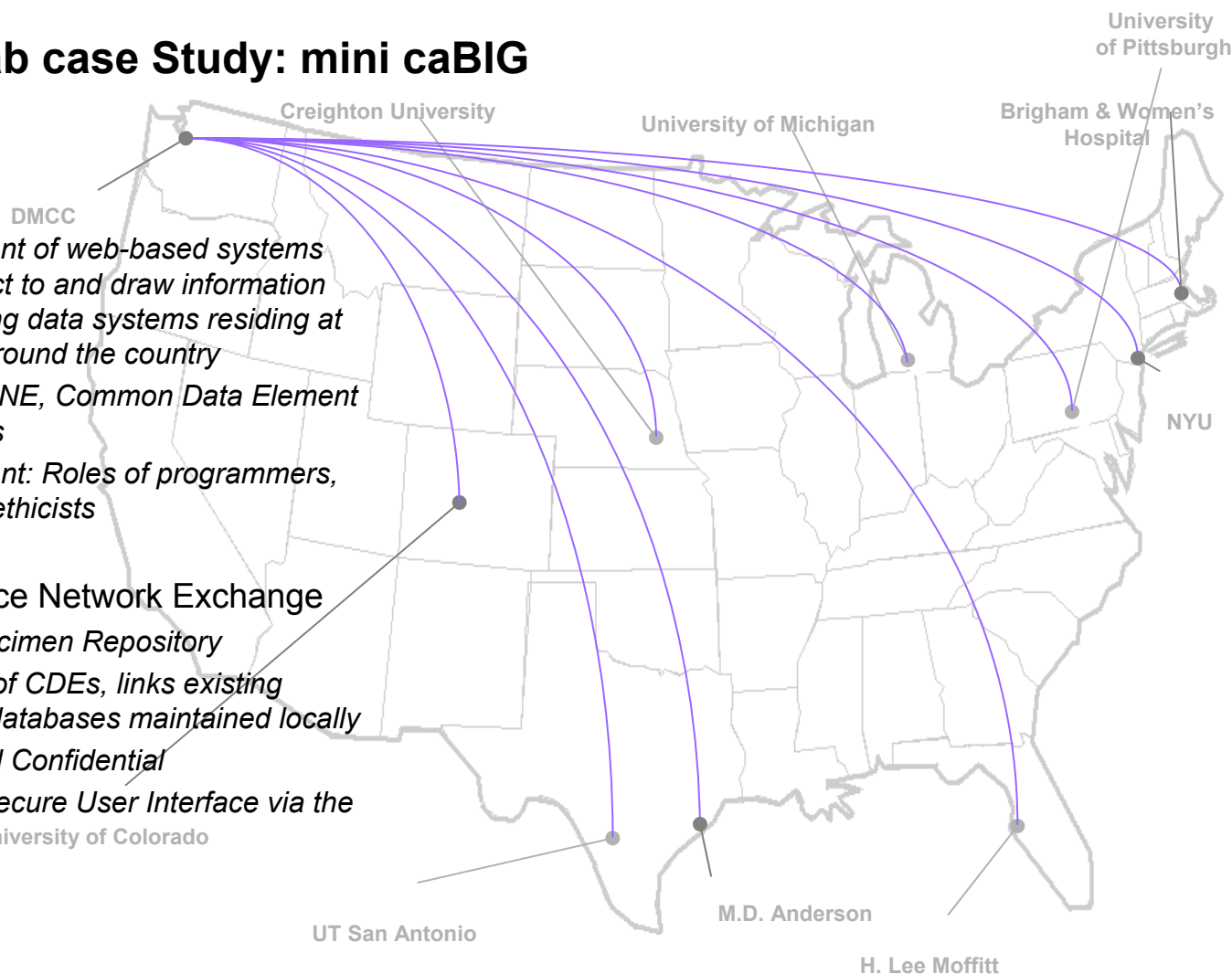
Interface

logic

data

Local Database

Export

External Database

Publication Control System

- Safe Export
- Export Schema
- Mapping tools
- Self Documenting Schema
- Exposed APIs
- Auxiliary Data Supported
- Identity Management Support
- Access Case Compliant

**caBIG Architecture Kickoff Meeting Presentation**
Fred Hutchinson Cancer Research Center

FRED
HUTCHINSON
CANCER
RESEARCH
CENTER
*Advancing Knowledge, Saving Lives*

# Thornquist Lab case Study: mini caBIG

University
of Pittsburgh

Creighton University

University of Michigan

Brigham & Women's
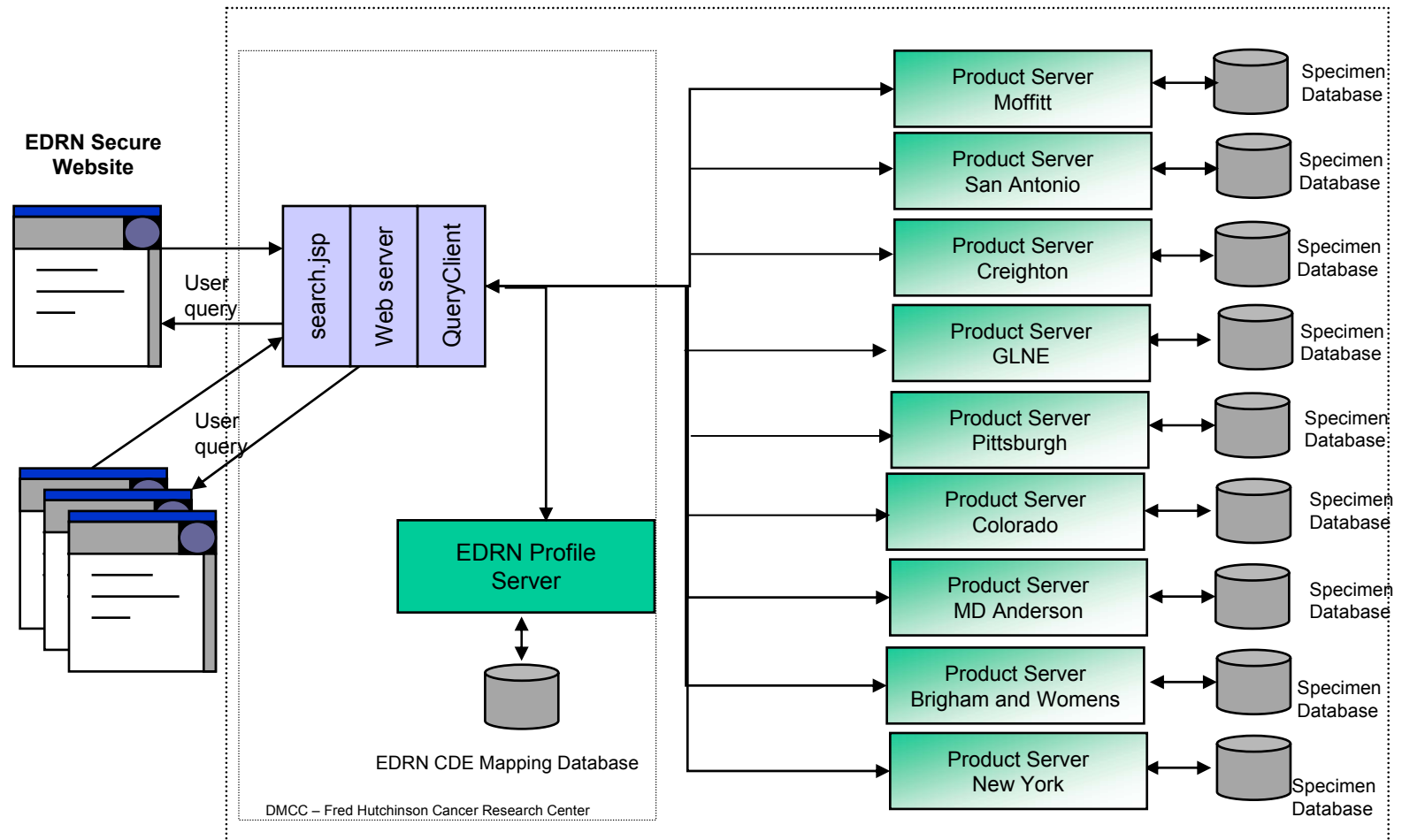Hospital

DMCC

- Structure
  - *Development of web-based systems that connect to and draw information from existing data systems residing at locations around the country*
  - *Pieces: ERNE, Common Data Element (CDE) tools*
  - *Development: Roles of programmers, scientists, ethicists*

NYU

- EDRN Resource Network Exchange
  - *Virtual Specimen Repository*
  - *By means of CDEs, links existing specimen databases maintained locally*
  - *Secure and Confidential*
  - *Dynamic Secure User Interface via the Internet*  University of Colorado

UT San Antonio

M.D. Anderson

H. Lee Moffitt

# Software Component Deployment

**caBIG Architecture Kickoff Meeting Presentation**
Fred Hutchinson Cancer Research Center

FRED
HUTCHINSON
CANCER
RESEARCH
CENTER
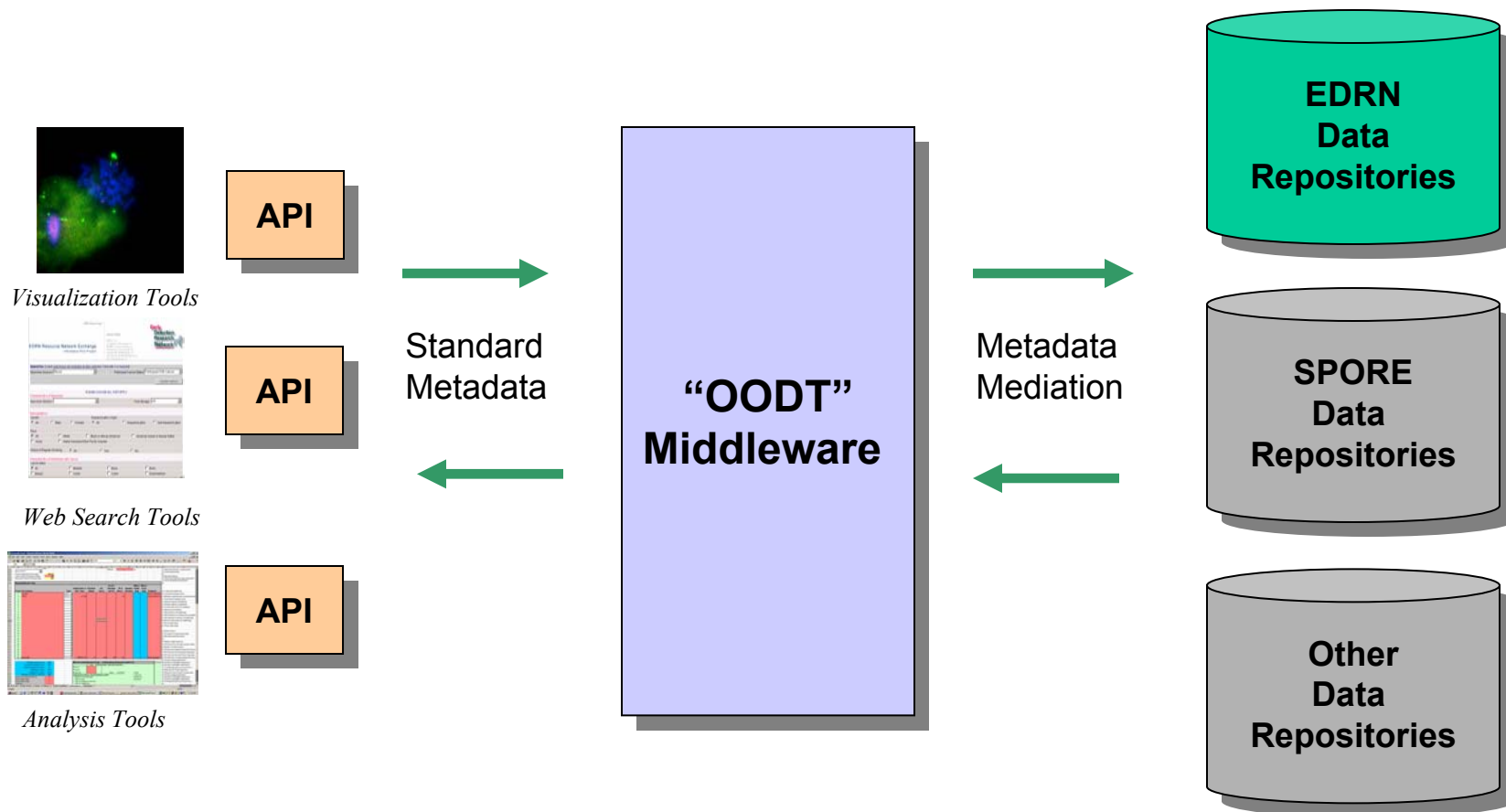*Advancing Knowledge, Saving Lives*

# EDRN Bioinformatics Architecture

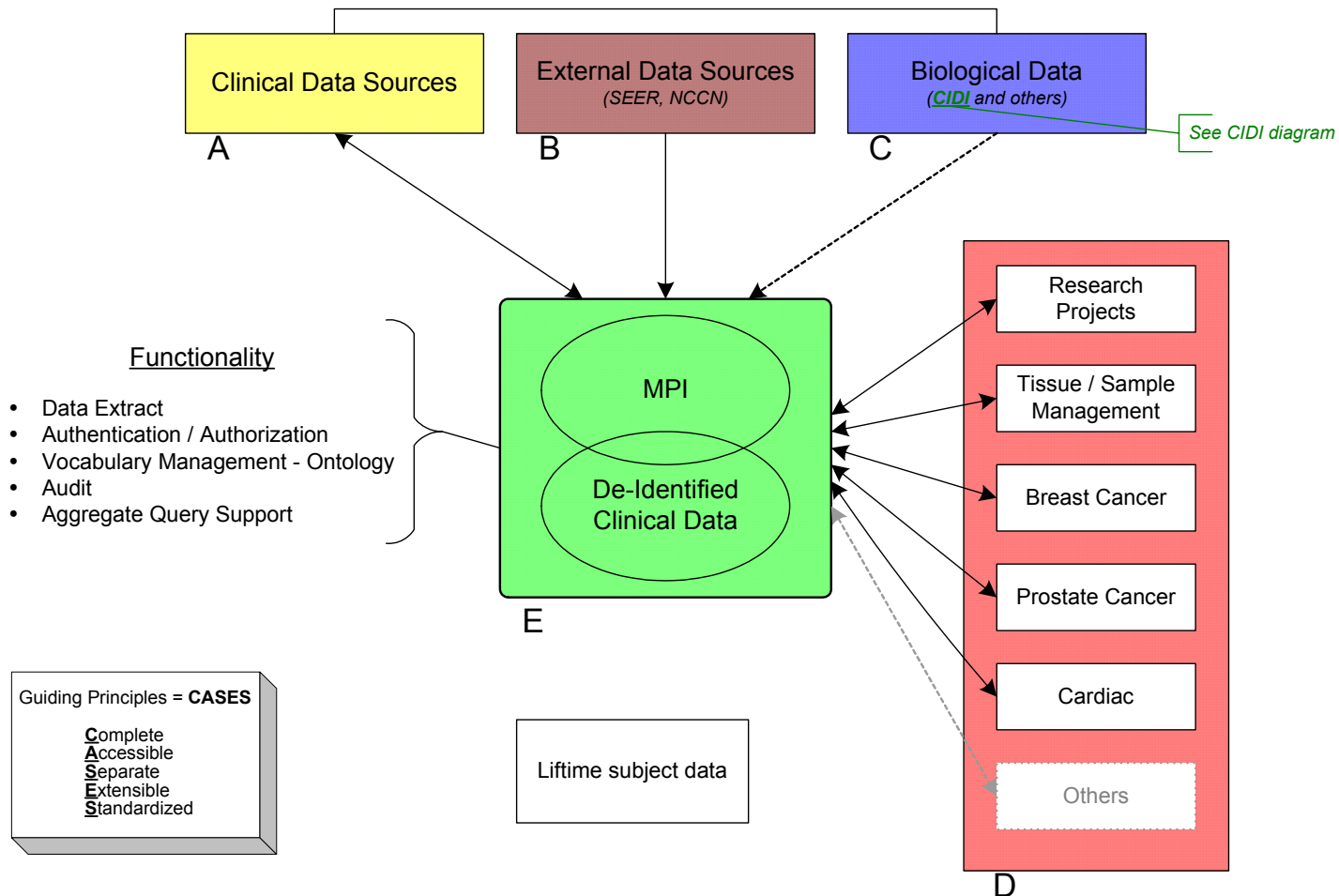**1. API's exposed** for Bioformatics tools and applications

**2. Middleware** creates the informatics infrastructure connecting systems and data

**3. Repositories** for storing and retrieving many data types data
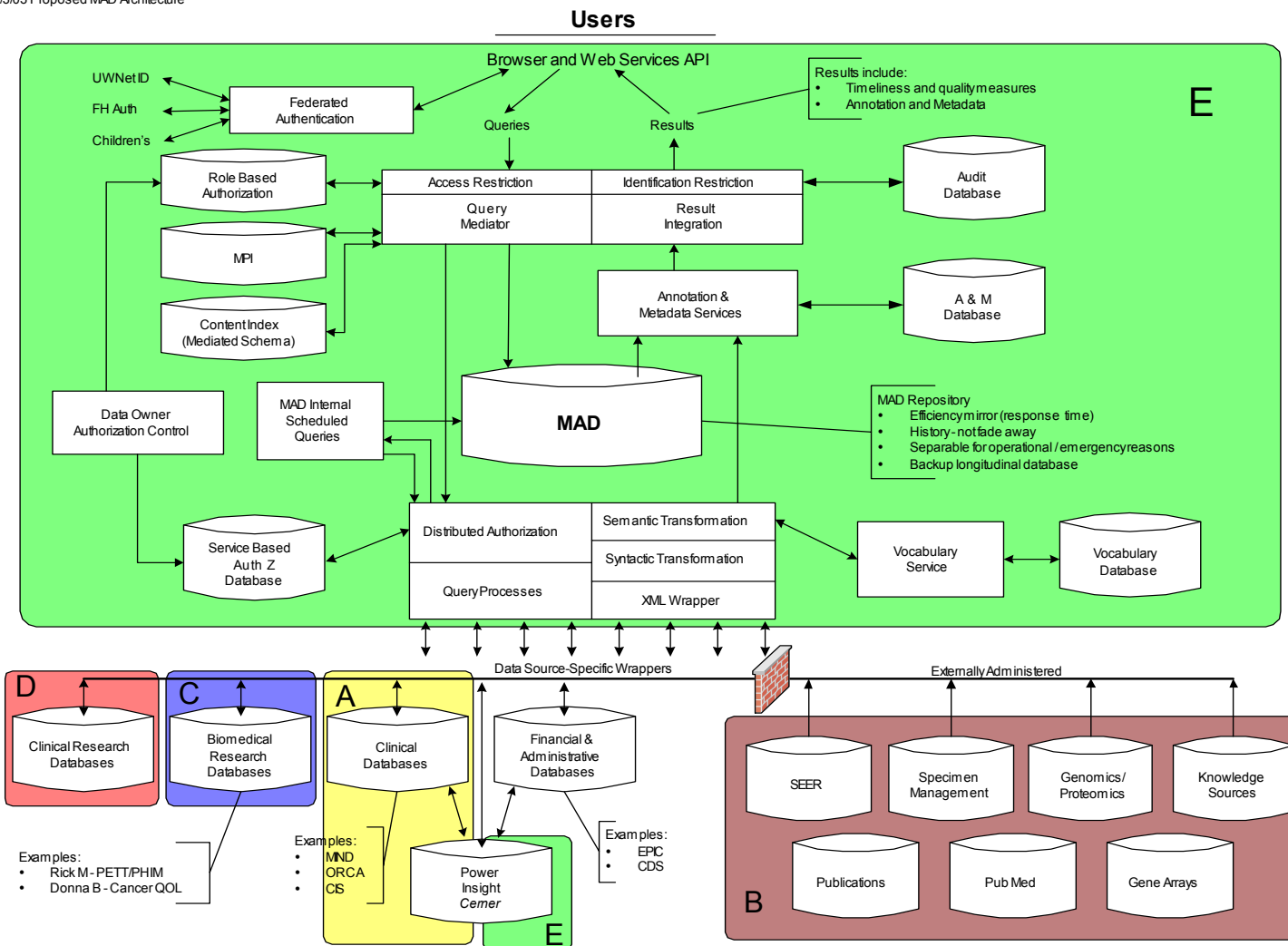


*Visualization Tools*

**API**

*Web Search Tools*

**API**

*Analysis Tools*

**API**

Standard Metadata

**"OODT" Middleware**

Metadata Mediation

**EDRN Data Repositories**

**SPORE Data Repositories**

**Other Data Repositories**

**caBIG Architecture Kickoff Meeting Presentation**
Fred Hutchinson Cancer Research Center

FRED
HUTCHINSON
CANCER
RESEARCH
CENTER
*Advancing Knowledge, Saving Lives*

# Reference Architecture – Conceptual Design

## High Level View

# caBIG Architecture Kickoff Meeting Presentation
Fred Hutchinson Cancer Research Center

FRED
HUTCHINSON
CANCER
RESEARCH
CENTER
*Advancing Knowledge, Saving Lives*

# Reference Architecture – Detailed View



4/3/03 Proposed MAD Architecture

# Summary

- Support the establishing and maintenance of common architecture
  - *Fostering alignment with a common vision in software design with an eye to collaboration*
  - *Development of tools that can interoperate between institutions/research initiatives*
  - *Understand the need to build and share these tools in a systematic way*

- Experience and Lessons Learned
  - *Managing and integrating systems from a variety of sources*
  - *Data publishing in real time as it becomes available*
  - *Challenges in supporting a variety of hardware and software systems*

- Flexibility is Essential
  - *Existing variability in data sets/systems/vocabularies/implementations that must be assembled in a grid environment*
  - *Depending on degree of expertise and budget available to the individual researcher*
  - *Based on the evolving nature of data elements in discovery oriented research*
  - *Based on the evolving nature of technology (connectivity, software platforms, hardware platforms)*